# Unsupervised Learning for Acoustic Applications
## Discovering Hidden Patterns in Sound

Can Evren Yarman

SLB, KTH

Nov 4, 2025

ASSA 2025

# What is Unsupervised Learning?

**Supervised Learning Algorithms:** Learn a mapping from input samples to output targets, commonly used for solving **classification** and **regression** problems.

$$data_{training} = \{input_i, label_i\} \Rightarrow Relationship$$

**Unsupervised Learning Algorithms:** Discover hidden patterns, structures, or groupings within data without relying on labeled outputs, commonly used for **clustering**, **dimensionality reduction**, and **representation learning** problems.

$$data_{training} = \{input_i\} \Rightarrow label_i$$

Association of these structures to patterns
$\leftrightarrow$ Choice of dictionary $\leftrightarrow$ Domain knowledge

# Question

Supervised vs Unsupervised:

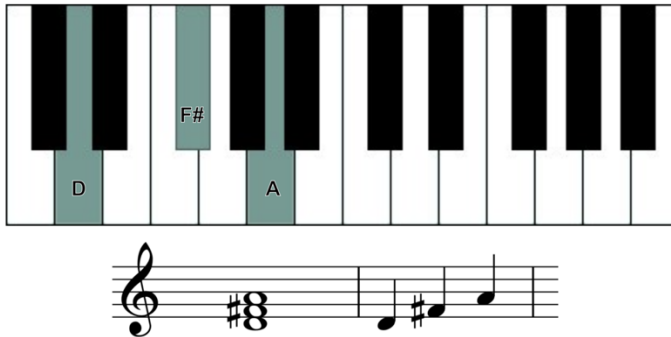In practice which one should precede the other? Why?

# Why Unsupervised Learning in Acoustics?

- Labeling acoustic data is expensive

- Useful for exploratory analysis

- Applications:
  - **Seismic interpretation:** Clustering groups seismic traces with similar waveforms to identify geological structures.
  - **Environmental sounds:** Clustering identifies distinct vocalization patterns in whale songs among marine mammals.
  - **Speech processing:** Clustering segments audio by speaker to distinguish different voices in conversations.
  - **Music software:** Clustering frequencies and durations for reverse engineering musical pieces
    - Spotify by Daniel Ek and Martin Lorentzon.
    - Auto-Tune by Andy Hildebrand, from Exxon Geologist to Autotune Inventor.

# Waveform Example

## Examples of acoustic waveforms from music

$$f(t) = \text{Re}\left\{\sum_m f_m e^{i\omega_m t}\right\}$$

## Time-Frequency Analysis of Musical Instruments*

Jeremy F. Alm[†]
James S. Walker[‡]

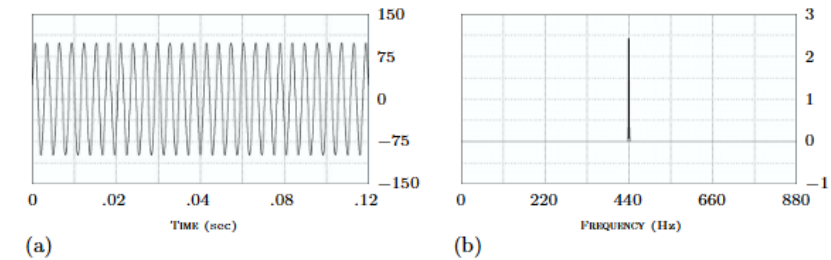**Fig. 2.1** *Fourier analysis of a pure tone. (a) Graph of a finite segment of a pure tone, 440 Hz. (b) Computer-calculated Fourier spectrum.*

(a) Piano note, $E_4$

(b) Spectrum of piano note

**Fig. 2.2** *Fourier analysis of the piano note $E_4$ (E above middle C). (Note: The vertical scales of all spectra shown in this paper have been normalized to the same range.)*

# Spectrogram Example

- Time-frequency representation (Importance of choice of dictionary)

$$f(t) = Re\left\{\sum_{m,n} f_{m,n} e^{i\omega_{mn}t} w(t - t_n)\right\}$$

- Useful for audio classification

- Example: song analysis



(a) Several piano notes

(b) Succession of windows

(c) Piano notes & single window

(d) Subsignal from single window

**Fig. 3.2** *Components of a spectrogram.*

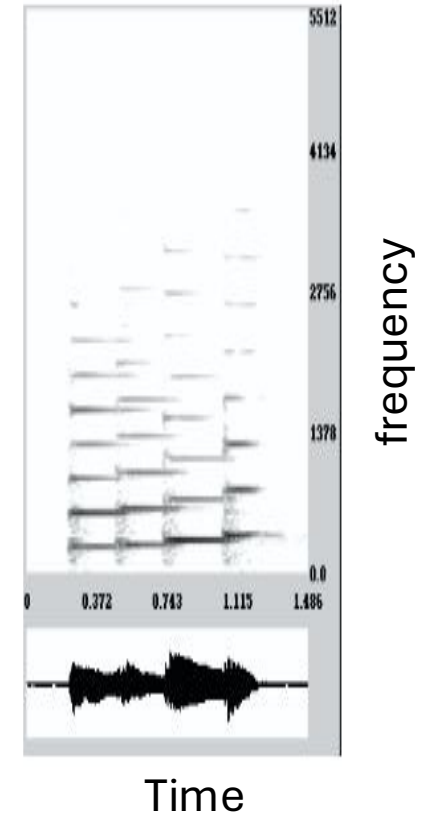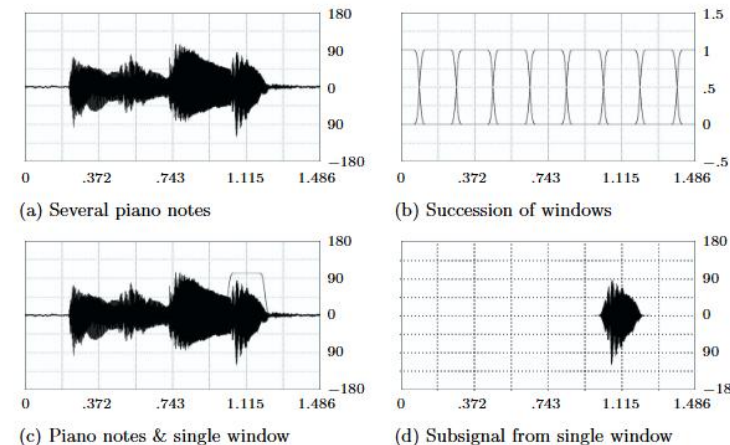# Types of Unsupervised Learning

- Clustering

- Dimensionality Reduction

- Association Rule Learning

# Hierarchical Clustering



## Description

- Group similar data points into clusters.
- There are two main types:
  - **Agglomerative (bottom-up)** – starts with each point as its own cluster and merges them.
  - **Divisive (top-down)** – starts with one big cluster and recursively splits it.

- Dendrograms for acoustic event grouping
- Example: Grouping similar acoustic emission events in material testing

## Algorithmic steps

1. Start with each data point as its own cluster (i.e., $n$ clusters).
2. Compute the distance (or similarity) between every pair of clusters using a linkage method.
3. Merge the two closest clusters based on the chosen linkage.
4. Update the distance matrix to reflect the merge.
5. Repeat steps 2–4 until all points are merged into a single cluster (or until a stopping criterion is met, like a desired number of clusters).
6. Cut the dendrogram at a chosen height.

# Hierarchical Clustering – Linkage methods

| Linkage Method | Distance Metric Used | Cluster Shape | Sensitivity to Outliers | Notes |
|---|---|---|---|---|
| Single Linkage | Minimum pairwise distance $D_{\min}(A, B)$ | Elongated | High | Can cause chaining |
| Complete Linkage | Maximum pairwise distance $D_{\max}(A, B)$ | Compact | High | Tends to break large clusters |
| Average Linkage | Average pairwise distance | Balanced | Moderate | Good general-purpose choice |
| Centroid | Distance between centroids | Elliptical | Moderate | May produce disconnected clusters |
| Ward's Method | Increase in within-cluster variance | Spherical | Low | Often preferred for numerical data |

$$D_{\min/\max}(A, B) = \min/\max\{d(a, b), a \in A, b \in B\} \qquad D_{\text{cent}}(A, B) = d(\mu_A, \mu_B)$$

$$D_{\text{avg}}(A, B) = |A|^{-1}|B|^{-1} \sum_{a \in A, b \in B} d(a, b) \qquad D_{\text{Ward}}(A, B) = |\boldsymbol{A}||\boldsymbol{B}|(|\boldsymbol{A} \cup \boldsymbol{B}|)^{-1} d(\mu_A, \mu_B) \quad \text{increase in within-cluster variance}$$

# Example – European cities



Example can be extended to earthquake locations using ObsPy.
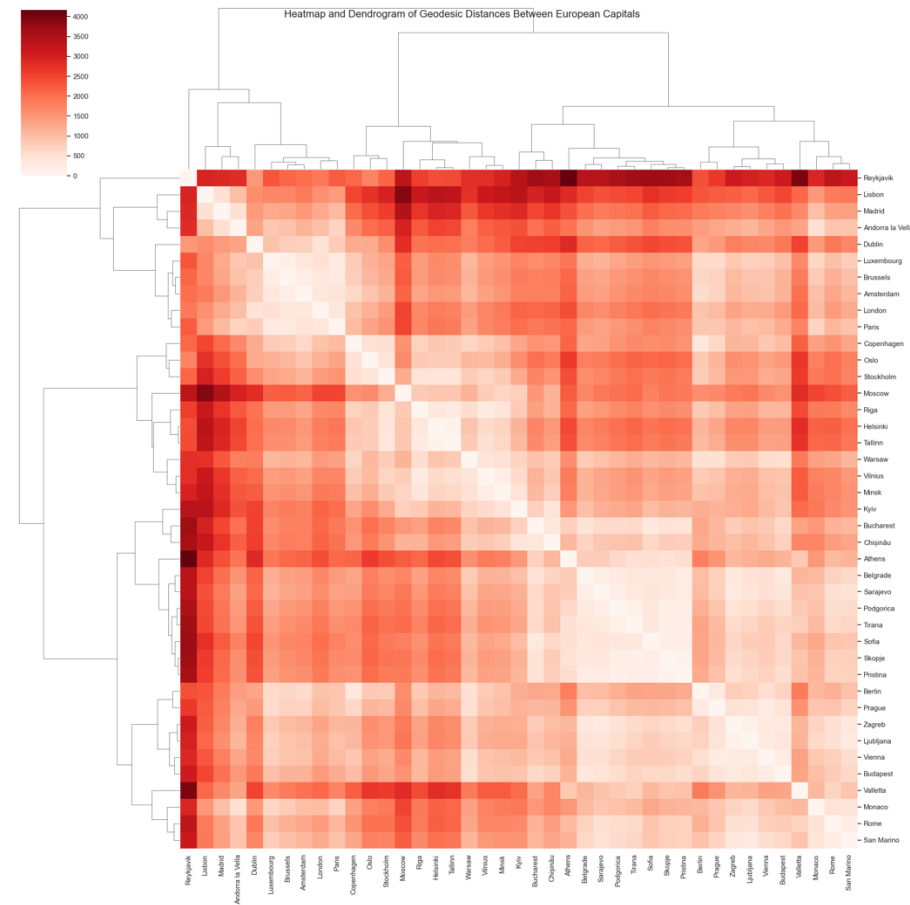
# Example – European cities (distances)

| | London | Paris | Brussels | Amsterdam | Berlin | Rome | Madrid | Lisbon | Vienna | Budapest | Prague | Warsaw | Copenhagen | Oslo | Stockholm | Helsinki | Tallinn | Riga | Vilnius | Athens | Sofia | Bucharest | Belgrade | Sarajevo | Zagreb | Ljubljana | Podgorica | Skopje | Tirana | Pristina | Reykjavik | Dublin | Moscow | Minsk | Kyiv | Chişinău | Valletta | Andorra la Vella | San Marino | Monaco | Luxembourg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| London | 0 | 343.92 | 321.6 | 358.97 | 934.52 | 1435.41 | 1263.1 | 1585.73 | 1238.91 | 1453.28 | 1037.54 | 1453.14 | 958.18 | 1156.06 | 1436.49 | 1825.86 | 1788.84 | 1681.46 | 1728.22 | 2395.06 | 2018.93 | 2096.72 | 1695.93 | 1624.51 | 1340.96 | 1231.29 | 1778.65 | 1945.51 | 1895.47 | 1881.96 | 1891.79 | 464.58 | 2508.57 | 1878.04 | 2140.05 | 2153.44 | 2089.05 | 1008.25 | 1261.64 | 1032.4 | 490.92 |
| Paris | 343.92 | 0 | 264.27 | 430.24 | 879.7 | 1106.6 | 1052.97 | 1454.35 | 1036.59 | 1247.67 | 885.46 | 1370.54 | 1028.58 | 1343.87 | 1546.63 | 1912.79 | 1862.97 | 1707.31 | 1700.65 | 2099.04 | 1761.37 | 1875.29 | 1449.04 | 1351.85 | 1082.42 | 966.59 | 1494.47 | 1669.62 | 1603.87 | 1611.68 | 2235.58 | 782.5 | 2493.65 | 1831.19 | 2029.65 | 1981.42 | 1748.74 | 708.59 | 948.8 | 689.87 | 287.76 |
| Brussels | 321.6 | 264.27 | 0 | 173.1 | 652.59 | 1173.91 | 1316.96 | 1714.25 | 917.34 | 1131.75 | 720.76 | 1163.4 | 766.54 | 1086.34 | 1282.8 | 1651.64 | 1603.56 | 1457.86 | 1469.31 | 2091.71 | 1701.41 | 1775.12 | 1377 | 1312.95 | 1026.33 | 919.69 | 1471.3 | 1633.78 | 1591.48 | 1568.41 | 2131.34 | 777.63 | 2260.06 | 1608.74 | 1841.8 | 1836.91 | 1849.94 | 952.15 | 981.23 | 823.93 | 187.38 |
| Amsterdam | 358.97 | 430.24 | 173.1 | 0 | 577.95 | 1296.48 | 1481.81 | 1865 | 937.72 | 1147.65 | 712.53 | 1096.73 | 622.27 | 914.67 | 1128 | 1505.27 | 1461.42 | 1334.33 | 1370.55 | 2165.19 | 1746.16 | 1791.39 | 1418.63 | 1377 | 1086.71 | 989.36 | 1542.12 | 1694.69 | 1666.95 | 1625.57 | 2021.12 | 759.09 | 2153.96 | 1519.24 | 1785.75 | 1818.6 | 1980.77 | 1125.18 | 1091.71 | 977.55 | 318.47 |
| Berlin | 934.52 | 879.7 | 652.59 | 577.95 | 0 | 1182.34 | 1871.43 | 2315.79 | 523.97 | 688.9 | 281.33 | 518.84 | 355.53 | 839.39 | 811.78 | 1107.6 | 1043.69 | 845.96 | 821.08 | 1803.12 | 1319.67 | 1297.02 | 1003.35 | 1032.18 | 768.86 | 723.22 | 1204.45 | 1316.01 | 1335.41 | 1240 | 2392.67 | 1321.27 | 1614 | 956.37 | 1208.16 | 1267.49 | 1848.86 | 1424.6 | 957.4 | 1072.24 | 603.4 |
| Rome | 1435.41 | 1106.6 | 1173.91 | 1296.48 | 1182.34 | 0 | 1367.61 | 1867.68 | 764.15 | 809.06 | 920.68 | 1315.68 | 1531.67 | 2007.22 | 1976.54 | 2202.87 | 2126.35 | 1867.47 | 1703.09 | 1052.78 | 896.4 | 1139.96 | 719.58 | 529.71 | 517.05 | 488.87 | 561.77 | 740.09 | 613.46 | 719.69 | 3305.16 | 1888.67 | 2379.56 | 1738.42 | 1678.27 | 1417.39 | 688.98 | 908.23 | 225.6 | 462.06 | 988.76 |
| Madrid | 1263.1 | 1052.97 | 1316.96 | 1481.81 | 1871.43 | 1367.61 | 0 | 503.7 | 1813.12 | 1978.45 | 1775.59 | 2293.99 | 2074.6 | 2389.87 | 2596.2 | 2952.8 | 2898.04 | 2717.34 | 2666.05 | 2375.28 | 2259 | 2478.93 | 2032.44 | 1861.99 | 1705.24 | 1601.64 | 1926.61 | 2107.31 | 1979.41 | 2082.53 | 2892.2 | 1450.64 | 3447.99 | 2770.39 | 2869.17 | 2703.23 | 1670 | 494.31 | 1388.04 | 991.02 | 1280.38 |
| Lisbon | 1585.73 | 1454.35 | 1714.25 | 1865 | 2315.79 | 1867.68 | 503.7 | 0 | 2303.35 | 2475.05 | 2247.71 | 2764.55 | 2480.71 | 2741.41 | 2992.99 | 3365.89 | 3317.16 | 3156.54 | 3126.12 | 2858.68 | 2761.18 | 2982.63 | 2535.99 | 2365.62 | 2206.15 | 2100.8 | 2428.08 | 2607.77 | 2477.42 | 2584.56 | 2950.79 | 1641.33 | 3914.85 | 3239.59 | 3359.71 | 3204.28 | 2114.13 | 994.6 | 1891.63 | 1493.31 | 1713.48 |
| Vienna | 1238.91 | 1036.59 | 917.34 | 937.72 | 523.97 | 764.15 | 1813.12 | 2303.35 | 0 | 214.59 | 251.14 | 556.25 | 870.82 | 1352.7 | 1242.7 | 1441.13 | 1363.69 | 1103.33 | 948.74 | 1282.56 | 818.14 | 858.22 | 492.32 | 508.78 | 267.72 | 277.98 | 680.69 | 796.06 | 811.59 | 721.05 | 2893.65 | 1687.09 | 1673.97 | 1006.02 | 1056.52 | 947.4 | 1375.69 | 1321.86 | 563.68 | 852.22 | 766.11 |
| Budapest | 1453.28 | 1247.67 | 1131.75 | 1147.65 | 688.9 | 809.06 | 1978.45 | 2475.05 | 214.59 | 0 | 443.19 | 545.01 | 1013.42 | 1483.6 | 1318.3 | 1461.36 | 1380.61 | 1106.43 | 910.63 | 1123.88 | 630.77 | 644.76 | 320.41 | 407.68 | 299.62 | 381.47 | 563.42 | 639.61 | 688.45 | 562.72 | 3079.02 | 1899.97 | 1572.74 | 931.16 | 901.34 | 745.06 | 1341.51 | 1484.29 | 647.88 | 996.34 | 980.12 |
| Prague | 1037.54 | 885.46 | 720.76 | 712.53 | 281.33 | 920.68 | 1775.59 | 2247.71 | 251.14 | 443.19 | 0 | 518.41 | 635.78 | 1119.91 | 1055.96 | 1304.64 | 1232.63 | 996.38 | 897.53 | 1533.39 | 1065.98 | 1081.56 | 742.06 | 754.41 | 487.53 | 446.86 | 927.07 | 1047.11 | 1058.07 | 972.19 | 2643.42 | 1470.48 | 1672.49 | 994.8 | 1145.48 | 1116.3 | 1574.84 | 1299.95 | 699.3 | 883.36 | 599.34 |
| Warsaw | 1453.14 | 1370.54 | 1163.4 | 1096.73 | 518.84 | 1315.68 | 2293.99 | 2764.55 | 556.25 | 545.01 | 518.41 | 0 | 673.39 | 1065.3 | 811.54 | 916.57 | 835.66 | 562.01 | 393.77 | 1597.26 | 1073.83 | 945.61 | 828.71 | 950.83 | 802.19 | 833.9 | 1097.39 | 1215.33 | 1063.69 |  | 2778 | 1832.85 | 1154.21 | 476.41 | 691.12 | 810.81 | 1885.79 | 1816.75 | 1119.9 | 1381.55 | 1083.92 |
| Copenhagen | 958.18 | 1028.58 | 766.54 | 622.27 | 355.53 | 1531.67 | 2074.6 | 2480.71 | 870.82 | 1013.42 | 635.78 | 673.39 | 0 | 484.13 | 523.35 | 885.26 | 789.42 | 727.01 | 816.04 | 2136.91 | 1638.15 | 1576.91 | 1333.4 | 1379.21 | 1122.77 | 1078.63 | 1550.35 | 1651.06 | 1680.89 | 1575.16 | 2111.24 | 1242.53 | 1565.75 | 1331.62 | 1483.03 |  | 2203.07 |  |  | 1377.65 | 1184.04 |
| Oslo | 1156.06 | 1343.87 | 1086.34 | 914.67 | 839.39 | 2007.22 | 2389.87 | 2741.41 | 1352.7 | 1483.6 | 1119.91 | 1065.3 | 484.13 | 0 | 417.79 | 789.42 | 789.22 | 845.18 | 1048.05 | 2606.43 | 2098.08 | 2007 | 1803.84 | 1859.9 | 1606.77 | 1561.91 | 2030.07 | 2123.21 | 2160.1 | 2046.32 | 1751.09 | 1268.59 | 1648.76 | 1218.11 | 1630.88 | 1860.62 | 2683.82 | 2036.45 | 1781.75 | 1813.72 | 1325.76 |
| Stockholm | 1436.49 | 1546.63 | 1282.8 | 1128 | 811.78 | 1976.54 | 2596.2 | 2992.99 | 1242.7 | 1318.3 | 1055.96 | 811.54 | 523.35 | 417.79 | 0 | 397.05 | 379.97 | 443.22 | 676.88 | 2407.91 | 1885.15 | 1744.96 | 1626.08 | 1721.66 | 1510.28 | 1495.92 | 1881.7 | 1941.87 | 2006.11 | 1866.28 | 2140.06 | 1633.3 | 1231.12 | 838.09 | 1266.9 | 1544.94 | 2617.72 | 2189.95 | 1755.03 | 1879.86 | 1325.76 |
| Helsinki | 1825.86 | 1912.79 | 1651.64 | 1505.27 | 1107.6 | 2202.87 | 2952.8 | 3365.89 | 1441.13 | 1461.36 | 1304.64 | 916.57 | 885.26 | 789.42 | 397.05 | 0 | 82.24 | 361.89 | 610.9 | 2468.74 | 1946.91 | 1753.45 | 1737.67 | 1867.33 | 1703.19 | 1714 | 2011.03 | 2035.73 | 2125.35 | 1964.48 | 2423.01 | 2029.84 | 895.06 | 715.9 | 1137.74 | 1486.81 | 2800.37 | 2525.88 | 1991.5 | 2173.86 | 1673.11 |
| Tallinn | 1788.84 | 1862.97 | 1603.56 | 1461.42 | 1043.69 | 2126.35 | 2898.04 | 3317.16 | 1363.69 | 1380.61 | 1232.63 | 835.66 | 789.42 | 789.22 | 379.97 | 82.24 | 0 | 279.66 | 529.9 | 2386.78 | 1864.81 | 1672.59 | 1655.72 | 1786.22 | 1624.76 | 1637.53 | 1929.32 | 1953.51 | 2043.38 | 1882.28 | 2456.84 | 2009.73 | 869.64 | 639.5 | 1065.61 | 1409.24 | 2720.28 | 2465.64 | 1916.23 | 2105.7 | 1617.67 |
| Riga | 1681.46 | 1707.31 | 1457.86 | 1334.33 | 845.96 | 1867.47 | 2717.34 | 3156.54 | 1103.33 | 1106.43 | 996.38 | 562.01 | 727.01 | 845.18 | 443.22 | 361.89 | 279.66 | 0 | 262.42 | 2108.79 | 1586.14 | 1400.03 | 1376.83 | 1510.32 | 1359.53 | 1380.07 | 1651.23 | 1673.85 | 1764.41 | 1602.68 | 2582.88 | 1959.73 | 844.43 | 403.68 | 837.61 | 1152.5 | 2447.7 | 2266.8 | 1662.68 | 1878.79 | 1443.6 |
| Vilnius | 1728.22 | 1700.65 | 1469.31 | 1370.55 | 821.08 | 1703.09 | 2666.05 | 3126.12 | 948.74 | 910.63 | 897.53 | 393.77 | 816.04 | 1048.05 | 676.88 | 610.9 | 529.9 | 262.42 | 0 | 1860.43 | 1340.9 | 1142.69 | 1154.23 | 1302.58 | 1186.77 | 1225.58 | 1432.27 | 1439.06 | 1539.16 | 1370.43 | 2798.66 | 2055.2 | 793.05 | 172.3 | 590.18 | 890.21 | 2246.24 | 2196.39 | 1511.86 | 1772.35 | 1421.29 |
| Athens | 2395.06 | 2099.04 | 2091.71 | 2165.19 | 1803.12 | 1052.78 | 2375.28 | 2858.68 | 1282.56 | 1123.88 | 1533.39 | 1597.26 | 2136.91 | 2606.43 | 2407.91 | 2468.74 | 2386.78 | 2108.79 | 1860.43 | 0 | 524.57 | 742.66 | 803.57 | 790.62 | 1080.38 | 1176.25 | 623.08 | 487.11 | 500.23 | 563.28 | 4167.72 | 2859.59 | 2231.98 | 1793.25 | 1486.3 | 1087.36 | 852.22 | 1948.8 | 1155.01 | 1294.85 | 1908.1 |
| Sofia | 2018.93 | 1761.37 | 1701.41 | 1746.16 | 1319.67 | 896.4 | 2259 | 2761.18 | 818.14 | 630.77 | 1065.98 | 1073.83 | 1638.15 | 2098.08 | 1885.15 | 1946.91 | 1864.81 | 1586.14 | 1340.9 | 524.57 | 0 | 295.56 | 327.72 | 418.66 | 680.46 | 794.96 | 334.88 | 174.52 | 317.62 | 176.78 | 3708.3 | 2479.03 | 1779.4 | 1284.17 | 1022.34 | 648.98 | 1069.72 | 1784.21 | 891.9 | 1294.85 | 1529.01 |
| Bucharest | 2096.72 | 1875.29 | 1775.12 | 1791.39 | 1297.02 | 1139.96 | 2478.93 | 2982.63 | 858.22 | 644.76 | 1081.56 | 945.61 | 1576.91 | 2007 | 1744.96 | 1753.45 | 1672.59 | 1400.03 | 1142.69 | 742.66 | 295.56 | 0 | 450.49 | 618.29 | 810.52 | 927.46 | 596.6 | 465.97 | 617.87 | 444.42 | 3682.28 | 2544.73 | 1501.27 | 1058.89 | 747.77 | 358.67 | 1365.27 | 1992.94 | 1091.81 | 1494.71 | 1618.04 |
| Belgrade | 1695.93 | 1449.04 | 1377 | 1418.63 | 1003.35 | 719.58 | 2032.44 | 2535.99 | 492.32 | 320.41 | 742.06 | 828.71 | 1333.4 | 1803.84 | 1626.08 | 1737.67 | 1655.72 | 1376.83 | 1154.23 | 803.57 | 327.72 | 450.49 | 0 | 192.5 | 368.47 | 485.96 | 278.83 | 319.74 | 387.69 | 242.92 | 3385.47 | 2154.32 | 1719.21 | 1136.36 | 983.42 | 697.81 | 1107.39 | 1544.29 | 644.41 | 1045.37 | 1207.56 |
| Sarajevo | 1624.51 | 1351.85 | 1312.95 | 1377 | 1032.18 | 529.71 | 1861.99 | 2365.62 | 508.78 | 407.68 | 754.41 | 950.83 | 1379.21 | 1859.9 | 1721.66 | 1867.33 | 1786.22 | 1510.32 | 1302.58 | 790.62 | 418.66 | 618.29 | 192.5 | 0 | 290.38 | 393.44 | 172.72 | 321.03 | 303.67 | 259.83 | 3317.58 | 2088.07 | 1901.69 | 1300.56 | 1205.35 | 959.1 | 943.99 | 1446.15 | 479.09 | 715.13 | 1133.77 |
| Zagreb | 1340.96 | 1082.42 | 1026.33 | 1086.71 | 768.86 | 517.05 | 1705.24 | 2206.15 | 267.72 | 299.62 | 487.53 | 802.19 | 1122.77 | 1606.77 | 1510.28 | 1703.19 | 1624.76 | 1359.53 | 1186.77 | 1080.38 | 680.46 | 810.52 | 368.47 | 290.38 | 0 | 117.59 | 458.44 | 608.89 | 586.9 | 542.14 | 3089.25 | 1803.14 | 1871.36 | 1222.46 | 1196.1 | 998.14 | 1107.98 | 1379.02 | 348.68 | 615.52 | 850.2 |
| Ljubljana | 1231.29 | 966.59 | 919.69 | 989.36 | 723.22 | 488.87 | 1601.64 | 2100.8 | 277.98 | 381.47 | 446.86 | 833.9 | 1078.63 | 1561.91 | 1495.92 | 1714 | 1637.53 | 1380.07 | 1225.58 | 1176.25 | 794.96 | 927.46 | 485.96 | 393.44 | 117.59 | 0 | 553.47 | 714.46 | 677.6 | 650.86 | 3001.46 | 1694.67 | 1936.55 | 1276.06 | 1282.57 | 1105 | 1128 | 1211.66 | 286.29 | 615.52 | 740.44 |
| Podgorica | 1778.65 | 1494.47 | 1471.3 | 1542.12 | 1204.45 | 561.77 | 1926.61 | 2428.08 | 680.69 | 563.42 | 927.07 | 1097.39 | 1550.35 | 2030.07 | 1881.7 | 2011.03 | 1929.32 | 1651.23 | 1432.27 | 623.08 | 334.88 | 596.6 | 278.83 | 172.72 | 458.44 | 553.47 | 0 | 185.18 | 131.01 | 158.68 | 3547.65 | 2243.08 | 1986.04 | 1414.7 | 1239.93 | 914.44 | 832.61 | 1456.15 | 578.03 | 973.73 | 1289.33 |
| Skopje | 1945.51 | 1669.62 | 1633.78 | 1694.69 | 1316.01 | 740.09 | 2107.31 | 2607.77 | 796.06 | 639.61 | 1047.11 | 1215.33 | 1651.06 | 2123.21 | 1941.87 | 2035.73 | 1953.51 | 1673.85 | 1439.06 | 487.11 | 174.52 | 465.97 | 319.74 | 321.03 | 608.89 | 714.46 | 185.18 | 0 | 153.13 | 76.89 | 3685.72 | 2409.1 | 1926.77 | 1399.18 | 1170.88 | 811.81 | 903.37 | 1639.74 | 762.72 | 1051.76 | 1454.8 |
| Tirana | 1895.47 | 1603.87 | 1591.48 | 1666.95 | 1335.41 | 613.46 | 1979.41 | 2477.42 | 811.59 | 688.45 | 1058.07 | 1063.69 | 1680.89 | 2160.1 | 2006.11 | 2125.35 | 2043.38 | 1764.41 | 1539.16 | 500.23 | 317.62 | 617.87 | 387.69 | 303.67 | 586.9 | 677.6 | 131.01 | 153.13 | 0 | 185.6 | 3675.85 | 2360.05 | 2061.36 | 1512.02 | 1308.17 | 959.1 | 758.91 | 1520.62 | 669.85 | 1051.76 | 1407.95 |
| Pristina | 1881.96 | 1611.68 | 1568.41 | 1625.57 | 1240 | 719.69 | 2082.53 | 2584.56 | 721.05 | 562.72 | 972.19 |  | 1575.16 | 2046.32 | 1866.28 | 1964.48 | 1882.28 | 1602.68 | 1370.43 | 563.28 | 176.78 | 444.42 | 242.92 | 259.83 | 542.14 | 650.86 | 158.68 | 76.89 | 185.6 | 0 | 3612.35 | 2344.83 | 1876.17 | 1335.61 | 635.9 | 776.51 | 944.35 | 1609.01 | 720.8 | 1121.86 | 1391.04 |
| Reykjavik | 1891.79 | 2235.58 | 2131.34 | 2021.12 | 2392.67 | 3305.16 | 2892.2 | 2950.79 | 2893.65 | 3079.02 | 2643.42 | 2778 | 2111.24 | 1751.09 | 2140.06 | 2423.01 | 2456.84 | 2582.88 | 2798.66 | 4167.72 | 3708.3 | 3682.28 | 3385.47 | 3317.58 | 3089.25 | 3001.46 | 3547.65 | 3685.72 | 3675.85 | 3612.35 | 0 | 1495.71 | 3317.19 | 2967.77 | 3381.47 | 3588.2 | 3975.45 | 2831.91 | 3109.54 | 2923.57 | 2317.35 |
| Dublin | 464.58 | 782.5 | 777.63 | 759.09 | 1321.27 | 1888.67 | 1450.64 | 1641.33 | 1687.09 | 1899.97 | 1470.48 | 1832.85 | 1242.53 | 1268.59 | 1633.3 | 2029.84 | 2009.73 | 1959.73 | 2055.2 | 2859.59 | 2479.03 | 2544.73 | 2154.32 | 2088.07 | 1803.14 | 1694.67 | 2243.08 | 2409.1 | 2360.05 | 2344.83 | 1495.71 | 0 | 2803.71 | 2217.13 | 2523.71 | 2576.87 | 2526.01 | 1336.54 | 1722.6 | 1465.94 | 954.36 |
| Moscow | 2508.57 | 2493.65 | 2260.06 | 2153.96 | 1614 | 2379.56 | 3447.99 | 3914.85 | 1673.97 | 1572.74 | 1672.49 | 1154.21 | 1565.75 | 1648.76 | 1231.12 | 895.06 | 869.64 | 844.43 | 793.05 | 2231.98 | 1779.4 | 1501.27 | 1719.21 | 1901.69 | 1871.36 | 1936.55 | 1986.04 | 1926.77 | 2061.36 | 1876.17 | 3317.19 | 2803.71 | 0 | 677.89 | 757.09 | 1146.06 | 2816.89 | 2970.36 | 2168.4 | 2216.84 | 2214.06 |
| Minsk | 1878.04 | 1831.19 | 1608.74 | 1519.24 | 956.37 | 1738.42 | 2770.39 | 3239.59 | 1006.02 | 931.16 | 994.8 | 476.41 | 1331.62 | 1218.11 | 838.09 | 715.9 | 639.5 | 403.68 | 172.3 | 1793.25 | 1284.17 | 1058.89 | 1136.36 | 1300.56 | 1222.46 | 1276.06 | 1414.7 | 1399.18 | 1512.02 | 1335.61 | 2967.77 | 2217.13 | 677.89 | 0 | 434.11 | 771.98 | 2241.87 | 2292.57 | 1560.63 | 1849.95 | 1547.6 |
| Kyiv | 2140.05 | 2029.65 | 1841.8 | 1785.75 | 1208.16 | 1678.27 | 2869.17 | 3359.71 | 1056.52 | 901.34 | 1145.48 | 691.12 | 1483.03 | 1630.88 | 1266.9 | 1137.74 | 1065.61 | 837.61 | 590.18 | 1486.3 | 1022.34 | 747.77 | 983.42 | 1205.35 | 1196.1 | 1282.57 | 1239.93 | 1170.88 | 1308.17 | 635.9 | 3381.47 | 2523.71 | 757.09 | 434.11 | 0 | 401.48 | 2066.13 | 2376.8 | 1543.82 | 1895.8 | 1742.26 |
| Chişinău | 2153.44 | 1981.42 | 1836.91 | 1818.6 | 1267.49 | 1417.39 | 2703.23 | 3204.28 | 947.4 | 745.06 | 1116.3 | 810.81 |  | 1860.62 | 1544.94 | 1486.81 | 1409.24 | 1152.5 | 890.21 | 1087.36 | 648.98 | 358.67 | 697.81 | 959.1 | 998.14 | 1105 | 914.44 | 811.81 | 959.1 | 776.51 | 3588.2 | 2576.87 | 1146.06 | 771.98 | 401.48 | 0 | 1715.16 | 2209.79 | 1325.52 | 1712.5 | 1703.98 |
| Valletta | 2089.05 | 1748.74 | 1849.94 | 1980.77 | 1848.86 | 688.98 | 1670 | 2114.13 | 1375.69 | 1341.51 | 1574.84 | 1885.79 | 2203.07 | 2683.82 | 2617.72 | 2800.37 | 2720.28 | 2447.7 | 2246.24 | 852.22 | 1069.72 | 1365.27 | 1107.39 | 943.99 | 1107.98 | 1128 | 832.61 | 903.37 | 758.91 | 944.35 | 3975.45 | 2526.01 | 2816.89 | 2241.87 | 2066.13 | 1715.16 | 0 | 1338.47 | 909.21 | 1060.2 | 1668.16 |
| Andorra la Vella | 1008.25 | 708.59 | 952.15 | 1125.18 | 1424.6 | 908.23 | 494.31 | 994.6 | 1321.86 | 1484.29 | 1299.95 | 1816.75 |  | 2036.45 | 2189.95 | 2525.88 | 2465.64 | 2266.8 | 2196.39 | 1948.8 | 1784.21 | 1992.94 | 1544.29 | 1446.15 | 1379.02 | 1211.66 | 1456.15 | 1639.74 | 1520.62 | 1609.01 | 2831.91 | 1336.54 | 2970.36 | 2292.57 | 2376.8 | 2209.79 | 1338.47 | 0 | 901.27 | 499.36 | 866.06 |
| San Marino | 1261.64 | 948.8 | 981.23 | 1091.71 | 957.4 | 225.6 | 1388.04 | 1891.63 | 563.68 | 647.88 | 699.3 | 1119.9 |  | 1781.75 | 1755.03 | 1991.5 | 1916.23 | 1662.68 | 1511.86 | 1155.01 | 891.9 | 1091.81 | 644.41 | 479.09 | 348.68 | 286.29 | 578.03 | 762.72 | 669.85 | 720.8 | 3109.54 | 1722.6 | 2168.4 | 1560.63 | 1543.82 | 1325.52 | 909.21 | 901.27 | 0 | 404.69 | 793.96 |
| Monaco | 1032.4 | 689.87 | 823.93 | 977.55 | 1072.24 | 462.06 | 991.02 | 1493.31 | 852.22 | 996.34 | 883.36 | 1381.55 | 1377.65 | 1813.72 | 1879.86 | 2173.86 | 2105.7 | 1878.79 | 1772.35 | 1294.85 | 1294.85 | 1494.71 | 1045.37 | 715.13 | 615.52 | 615.52 | 973.73 | 1051.76 | 1051.76 | 1121.86 | 2923.57 | 1465.94 | 2216.84 | 1849.95 | 1895.8 | 1712.5 | 1060.2 | 499.36 | 404.69 | 0 | 660.32 |
| Luxembourg | 490.92 | 287.76 | 187.38 | 318.47 | 603.4 | 988.76 | 1280.38 | 1713.48 | 766.11 | 980.12 | 599.34 | 1083.92 | 1184.04 | 1325.76 | 1325.76 | 1673.11 | 1617.67 | 1443.6 | 1421.29 | 1908.1 | 1529.01 | 1618.04 | 1207.56 | 1133.77 | 850.2 | 740.44 | 1289.33 | 1454.8 | 1407.95 | 1391.04 | 2317.35 | 954.36 | 2214.06 | 1547.6 | 1742.26 | 1703.98 | 1668.16 | 866.06 | 793.96 | 660.32 | 0 |

# Example – European cities (dendrogram)



Agglomerative (bottom-up)

Hierarchical Clustering of European Capitals (Geodesic Distances)

Divisive (top-down)

# Example – European cities

## (heat map + dendrogram)



Heatmap and Dendrogram of Geodesic Distances Between European Capitals

# Example – European cities

**k=4**



Hierarchical Clustering of European Capitals (4 Clusters)

**k=29**



Hierarchical Clustering of European Capitals (29 Clusters)

# Hierarchical Clustering - Properties

- Deterministic

- Computational cost $\mathcal{O}(n^3)$

- Dependent on distance definition

- Not necessarily optimal!
  Q: Why?
  A: Greedy at each iteration

# K-Means Clustering

## Description

- Partitions data into k-clusters

$$\arg\min_{\mu_k} \sum_{k=1}^{K} \sum_{x_i}^{\Omega_k} D(x_i, \mu_k)$$

- $D(x_i, \mu_k) = [x_i - \mu_k]^T [x_i - \mu_k]$
- $D(x_i, \mu_k) \neq$ **variance ! Q: Implications?**

- Assumes spherical clusters
- Fast and simple

## Algorithmic Steps

1. Choose the number of clusters (K) you want to find.

2. Initialize K centroids $\mu_k$ **randomly** (these are the centers of your clusters).

3. Assign each data point to the nearest centroid (based on distance, usually Euclidean): $x_i \in \Omega_k$

4. Update centroids by calculating the mean of all points assigned to each cluster: $\mu_k = |\Omega_k|^{-1} \sum_{x_i \in \Omega_k} x_i$

5. Repeat steps 3 and 4 until the centroids no longer change significantly (i.e., convergence).

# Example – European cities

# Example – European cities

# Example – European cities

# Example – European cities

# Example – European cities

# Example – European cities



K-Means Clustering of European Capitals (k=4) with Final Centroids

# K-Means vs Hierarchical clustering (k=4)



K-Means Clustering of European Capitals (4 Clusters)

Hierarchical Clustering of European Capitals (4 Clusters)

# K-mean vs Hierarchical clustering (k=29)



K-Means Clustering of European Capitals (29 Clusters)      Hierarchical Clustering of European Capitals (29 Clusters)

# K-means clustering - properties

- Non-deterministic
  - Choice of initial centroids
    - Informed decision / corners / sufficiently distant / multiple realization
  - Choice of # of clusters
    - From application / SME

- Depends on the distance definition: Euclidean

- Similar to Ward's method, which minimizes the increase in total within-cluster variance (similar to k-means). Weighting makes the difference: penalizes larger clusters with large variations.

# Gaussian Mixture Model (GMM) & Expectation Maximization (EM)

## GMM Description

- Probabilistic clustering minimizing log likelihood:

$$\log L = \sum_{i=1}^{n} \log\left( \sum_{k=1}^{K} \pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k) \right)$$

$$\mathcal{N}(x_i | \mu_k, \Sigma_k) = \left( (2\pi)^n \det \Sigma_k \right)^{-1/2} e^{(x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k)/2}$$

- Soft assignments
- Handles elliptical clusters
- Uses Expectation-Maximization

## EM Algorithmic steps

- E-Step (Expectation):
  Compute the **responsibility** $\gamma_{ik}$ —the probability that point $x_i$ belongs to cluster $k$

  $$\gamma_{ik} = \frac{\pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k)}{\sum_l \pi_l \mathcal{N}(x_i | \mu_l, \Sigma_l)}$$

- M-Step (Maximization):
  1. Effective number of points assigned to cluster: $N_k = \sum_{i=1}^{n} \gamma_{ik}$
  2. Update means: $\mu_k = N_k^{-1} \sum_{i=1}^{n} \gamma_{ik} x_i$
  3. Update covariance:
     $$\Sigma_k = N_k^{-1} \sum_{i=1}^{n} \gamma_{ik} (x_i - \mu_k)(x_i - \mu_k)^T$$
  4. Update mixing coefficient: $\pi_k = N_k/n$

- Repeat EM steps until $\log L$ is below a threshold.

# Example – European cities



GMM Clustering of European Cities (k=4)

# GMM clustering + EM - properties

- EM is non-deterministic
  - Initial centroids & (inv) covariance matrices
    - Informed decision / corners / sufficiently distant / multiple realization
    - Modes, maxima, etc. of histogram
  - Choice of optimization scheme
    - Greedy + Stochastic descent
  - # of clusters
    - From application / SME
  - Assumptions – sufficiently distinct peaks
- Alternative
  - Histogram decomposition
  - Challenge handling multivariate histograms

Gaussian mixture model decomposition of multivariate signal

Gustav Zickert[1] · Can Evren Yarman[2]

Signal, Image and Video Processing (2022) 16:429–436
https://doi.org/10.1007/s11760-021-01961-y

# DBSCAN

## Description

- Groups together points that are closely packed and marks points that lie alone in low-density regions as outliers.
  ~ # of modes in histogram

- Does not require specifying the number of clusters in advance.

- Can find arbitrarily shaped clusters.

- Can identify noise/outliers.

- Sensitive to the choice of ε and minPts.

- Struggles with clusters of varying density.

- Example: Detecting anomalous acoustic events (e.g., microseismic outliers)

## Algorithmic steps

1. For each point in the dataset:
   - Find all points within distance ε (its neighborhood).
   - If the number of neighbors ≥ minPts, mark it as a core point.

2. Form a cluster:
   - For each core point not yet assigned to a cluster:
     - Create a new cluster.
     - Add the core point and all its density-reachable points (i.e., points within ε of the core or other reachable core points).

3. Repeat until all points are visited.

4. Label remaining points that are not part of any cluster as noise.

# Example – European cities



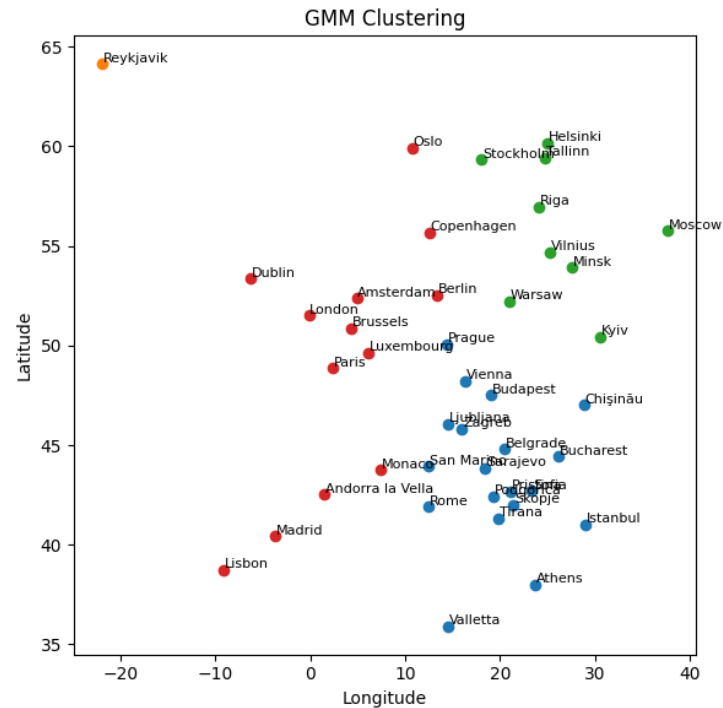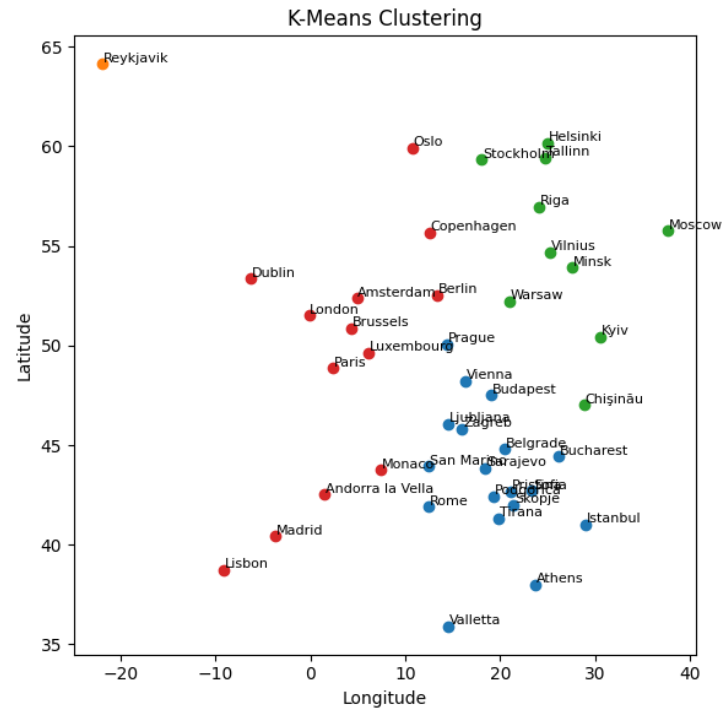DBSCAN Clustering of European Cities

# DBScan Choice of parameters

Convert degrees to kilometers:

- 1° latitude ≈ 111 km.

- 1° longitude ≈ 111 km × cos(latitude).

- So eps = 0.5° ≈ 55 km — too small for cities spread across Europe.

- For clustering cities into regional groups, you need eps in the range of 3° to 6° (≈ 300–600 km).

# Comparison of clustering methods

# Choosing the Right Clustering Algorithm

| Feature | K-Means | GMM (Gaussian Mixture Models) | DBSCAN | Hierarchical Clustering |
|---|---|---|---|---|
| **Cluster Shape Assumption** | Spherical, equal size | Elliptical, variable size | Arbitrary shapes | Arbitrary shapes |
| **Cluster Assignment** | Hard (each point belongs to one) | Soft (probabilistic) | Hard | Hard |
| **Number of Clusters Required** | Yes | Yes | No | Optional (can cut dendrogram) |
| **Noise Handling** | Poor | Poor | Good (identifies noise points) | Poor |
| **Scalability** | High | Moderate | Moderate | Low (especially with large datasets) |
| **Interpretability** | Easy | Moderate | Moderate | Easy (via dendrograms) |
| **Distance Metric** | Euclidean | Mahalanobis (can vary) | Any (usually Euclidean) | Any (usually Euclidean) |
| **Best Use Case in Acoustics** | Speaker clustering, waveform grouping | Speaker modeling, seismic facies | Anomaly detection in acoustic signals | Grouping similar acoustic events |
| **Strengths** | Fast, simple | Flexible, probabilistic | Detects noise, no need for k | Reveals hierarchy, no need for k |
| **Limitations** | Sensitive to initial centroids | Computationally intensive | Sensitive to parameters (ε, minPts) | Computationally expensive |
| **Computational complexity** | $\mathcal{O}(nkdi)$ | $\mathcal{O}(nkd^2i)$ | $\mathcal{O}(n\log n)$ [1] to $\mathcal{O}(n^2)$ | $\mathcal{O}(n^2)$ [2] to $\mathcal{O}(n^3)$ |

- **GMM** → Best for modeling complex distributions.
- **K-Means** → Acceptable for moderate dimensions.
- **Hierarchical** → Not ideal.
- **DBSCAN** → Worst in high dimensions.

$n$ = number of points
$k$ = number of clusters
$d$ = dimensionality
$i$ = number of iterations

[1] with spatial indexing like KD-tree or R-tree.
[2] with efficient data structures (e.g., using priority queues).

# Evaluation Metrics

| Metric | Formula | Type | Best Value | Use Case |
|--------|---------|------|-----------|----------|
| Silhouette Score | $$s(i) = \frac{b(i) - a(i)}{\max\big(a(i), b(i)\big)}$$ | Internal | Close to 1 | General clustering quality |
| Davies-Bouldin Index | $$DBI = \frac{1}{k}\sum_{i=1}^{k}\max_{j \neq i}\left(\frac{\sigma_i - \sigma_j}{d\big(\mu_i, \mu_j\big)}\right)$$ | Internal | Close to 0 | Compactness & separation |
| AIC / BIC | $$\text{AIC} = 2k - 2\log(L)$$ $$\text{BIC} = k\log(n) - 2\log(L)$$ | Model-based | Lower is better | GMM model selection |

$a(i)$ = average distance from point $i$ to all other points in the same cluster.
$b(i)$ = minimum average distance from point $i$ to points in other clusters.
$\sigma(i)$ = average distance of all points in cluster $i$ to its centroid $\mu_i$.

- **Akaike Information Criterion**: $\text{AIC} = 2k - 2\log(L)$
- **Bayesian Information Criterion**: $\text{BIC} = k\log(n) - 2\log(L)$

$k$: # of model parameters, $n$: # of data points, $L$: likelihood function
Evaluates how well a probabilistic model (like GMM) fits the data, while penalizing model complexity.

# Real-World Applications

- Seismic characterization
  - https://www.marine-geo.org/doi/10.60521/331931
  - https://wiki.seg.org/wiki/Open_data
- Speaker diarization
  - https://dihardchallenge.github.io/dihard3/
  - https://github.com/wq2012/awesome-diarization
- Environmental sound clustering
  - **The Marinexplore and Cornell University Whale Detection Challenge** https://www.kaggle.com/competitions/whale-detection-challenge/data

# Challenges and Future Directions

- Lack of ground truth

- Interpretability

- Evaluation without labels

- Self-supervised learning

- High-dimensionality of audio features
  - Temporal & directional dependencies

# Recent & Evolving Trends

- Self-supervised learning in audio (e.g., wav2vec)
- Integration with deep learning (autoencoders, contrastive learning)

# Summary and Q&A

- Unsupervised learning reveals **hidden** (expected or unexpected) patterns

  !!! Hidden ≠ Desired or Useful !!!

- Key techniques:
  - choice of representation (i.e. features)
  - dimension reductions (relevant features)
  - clustering (of features)

- Applications in seismic, speech, and audio

- Question the status quo

# Tutorial – Suggestion (others welcomed)

- **Load audio signals** (whale sounds).
- **Listen the audio signal**
- **Computes time varying features (i.e. STFT, Wavelet transform)**
- Extracts **feature vectors** by taking the mean of the feature across time.
  - Other features (e.g., mean, variance, or log-mel spectrogram).
- **Standardizes features** for better clustering performance.
- **Reduces dimensionality** with PCA (2 components for visualization).
- **Clusters signals** using Kmeans, etc.
- **Plots clusters** in PCA space.
- **V&V**
  - **Verification:** Are we building solution right
  - **Validation:** Are we building the right solution